

Convolutional neural network-based pelvic floor structure segmentation using magnetic resonance imaging in pelvic organ prolapse

Fei Feng

University of Michigan-Shanghai Jiao Tong University Joint Institute, Shanghai Jiao Tong University, Shanghai 200240, China

James A. Ashton-Miller

Department of Mechanical Engineering, University of Michigan, Ann Arbor, MI 48109, USA

John O. L. DeLancey

Department of Obstetrics and Gynecology, University of Michigan, Ann Arbor, MI 48109, USA

Jiajia Luo^{a)}

Biomedical Engineering Department, Peking University, Beijing 100191, China

(Received 30 December 2019; revised 18 June 2020; accepted for publication 22 June 2020; published 28 July 2020)

Purpose: Automated segmentation could improve the efficiency of modeling-based pelvic organ prolapse (POP) evaluations. However, segmentation performance is limited by the blurry soft tissue boundaries. In this study, we aimed to present a hybrid solution for uterus, rectum, bladder, and levator ani muscle segmentation by combining a convolutional neural network (CNN) and a level set method.

Methods: We used 24 sagittal pelvic floor magnetic resonance (MR) series from six anterior vaginal prolapse and six posterior vaginal prolapse subjects (a total 528 MR images). The stress MR images were performed both at rest and at maximal Valsalva. We assigned 264 images for training, 132 images for validation, and 132 images for testing. A CNN was designed by introducing a multi-resolution features pyramid module (MRFP) into an encoder-decoder model. Depth separable convolution and pretraining were used to improve model convergence. Multiclass cross entropy loss and multiclass Dice loss were used for model training. The dice similarity coefficient (DSC) and average surface distance (ASD) were used for evaluating the segmentation results. To prove the effectiveness of our model, we compared it with advanced segmentation methods including Deeplabv3+, U-Net, and FCN-8s. The ablation study was designed to quantify the contributions of MRFP, the encoder network, and pretraining. Besides, we investigated the working mechanism of MRFP in the segmentation network by comparing our model with three of its variants. Finally, the level set method was used to improve the CNN model further.

Results: Dice loss showed better segmentation performance than multiclass cross entropy loss. MRFP was efficacious for different encoder networks. With MRFP, U-Net and U-Net-X (X represents Xception encoder network) have improved the DSC, on average by 6.8 and 5.3 points. Compared with different CNN models, our model achieved the highest average DSC of 65.6 points and the lowest average ASD of 2.9 mm. With the level set method, the DSC of our model improved to 69.4 points.

Conclusions: MRFP proved to be effective in addressing the blurry soft tissue boundary problem on pelvic floor MR images. A hybrid solution based on CNN and level set method was presented for pelvic organ segmentation both at rest and at maximal Valsalva; with this method, we achieved state-of-the-art results. © 2020 American Association of Physicists in Medicine [<https://doi.org/10.1002/mp.14377>]

Key words: convolutional neural network, level set, MRI segmentation, multiresolution features pyramid, pelvic organ prolapse

1. INTRODUCTION

Pelvic organ prolapse (POP) is an abnormal caudal displacement and deformation of one or more female pelvic floor organs. Pelvic organ prolapse can cause considerable discomfort to women both physically and mentally. In the United States, about 200 000 women undergo POP surgery every year, at a total cost of more than \$1 billion.^{1,2} The most common imaging techniques to evaluate POP include magnetic resonance (MR) and ultrasound imaging. Due to the good contrast of soft tissues, MR imaging has always been the golden

standard for organ segmentation. Organ segmentation is crucial for three-dimensional (3D) geometric model reconstruction, finite element simulation of POP, and surgical planning.^{3,4} Currently, manual organ segmentation is still the most widely used technique. However, the manual segmentation is not only time-consuming but also susceptible to large inconsistencies depending on the experience and skill of the evaluators and the quality of MR scans. To speed up the segmentation process, computer-aided diagnostic techniques may hold promise.

Several difficulties constrain the pelvic organ segmentation performance. First, MR images do not provide high

enough contrast at the boundary of each organ, which makes segmentation particularly challenging for humans. Second, the occurrence rate is unbalanced between organs, which limits model convergence. For example, organs like the bladder are present in more MR images, whereas some organs, including rectum and uterus, may not be seen at all in many MR images, when viewed laterally. Adding to that challenge, some patients have undergone hysterectomy and lack a uterus. Third, large variations exist in these data. For instance, the shape and size of pelvic organs vary widely between resting and stressed (Valsalva) states (Fig. 1). Besides, the levator ani muscle exhibits a large inter-subject variance on MR images due to its structural complexity.

Computer-aided segmentation techniques include both deep learning and non-deep learning methods. The non-deep learning methods, including the deformable model and level set methods, have played an important role in the segmentation of the cardiac ventricle and other human body regions.^{5–8} One limitation is that those methods often fail to converge for images with blurry boundaries. Besides, their segmentation speeds do not fulfill the current needs for rapid segmentation as they require much human interaction. Moreover, the poor generalization is a typical problem that both automatic and semi-automatic methods face. Generalization problems are usually related to generalization in new regions or on new data. The first generalization problem means that one organ segmentation algorithm is usually not suitable for another organ. This hampers POP analysis since we usually want to obtain a segmentation of the uterus, bladder, levator ani muscle, rectum, vaginal walls, and other tissues simultaneously. The second generalization problem is even more crucial for the clinical application of automatic segmentation tools.

Since there are large variations in the structural profiles, it is challenging to find a solution that can adapt to inter-subject variability in MR images.

Recently, the convolutional neural network (CNN) has become the mainstream method for approaching many computer vision and medical imaging analysis problems. These include cell, lesion, tumor, retinal vessel, cardiac structure, and brain segmentation.^{9–13} Compared with non-deep learning methods, CNN usually does not rely on much prior knowledge of the data,^{14,15} and it is trained with MR data from different subjects. Thus it has good generalization performance. The basic idea of the CNN method is that it uses several convolution layers to extract features so it can provide pixelwise segmentation. Some researchers have proved that the sequentially stacked convolution layers are difficult to converge, so the residual connection and shortcut connection were proposed in ResNet¹⁶ and U-Net¹⁰ respectively, to smooth the model training process and preserve more detailed information.

Several CNN models were designed for different segmentation problems. U-Net¹⁰ adopted the encoder-decoder network to accomplish neuronal structures segmentation and cell tracking tasks. V-net¹⁷ used a 3D convolution to accomplish the volumetric segmentation task. DeepMedic¹¹ employed a dual-path 3D CNN based on dense patch ideas to deal with the high computational burden when training 3D CNN for brain lesion segmentation. UNet++¹⁸ connected the encoder and decoder networks by a series of dense skip connections to avoid eliminating the gap between encoder and decoder networks and obtained better performance than U-Net and wide U-Net on four segmentation datasets.

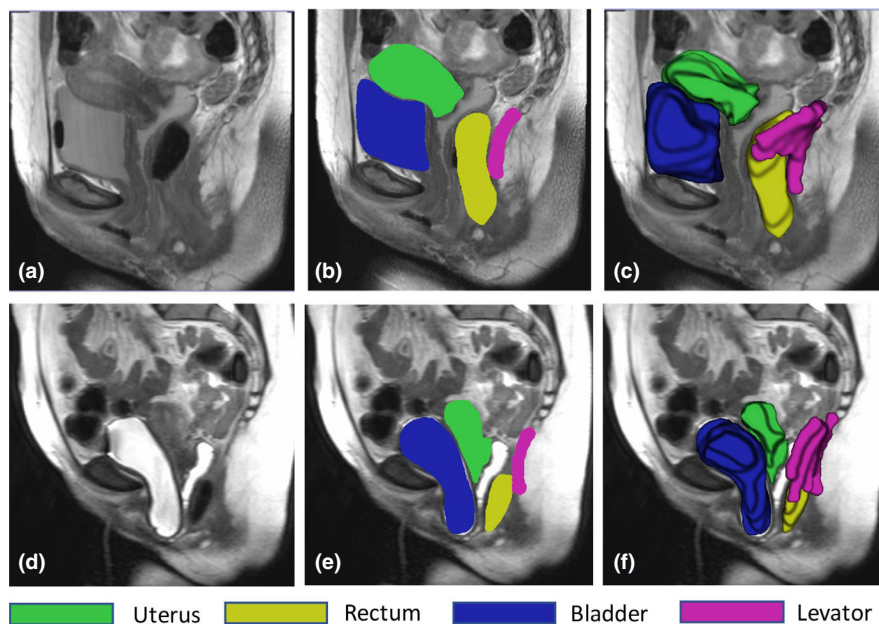


FIG. 1. Left lateral views of a patient with anterior vaginal wall prolapse. (a) and (d) Midsagittal magnetic resonance images at rest and at maximum Valsalva. (b) and (e) Similar images of the pelvic floor organs, including the uterus, rectum, bladder, and levator ani muscle, shown at rest and at maximum Valsalva. (c) and (f) Views of the three-dimensional models of the pelvic floor organs. [Color figure can be viewed at wileyonlinelibrary.com]

However, these designs could not capture different scales of semantic information. Segmentation is a task that needs details at different scales. Coarse segmentation could be achieved from lower resolution feature maps, while the fine-grained boundary information must be detected from higher resolution feature maps. Therefore, different sizes of features may preserve different scales of context information.¹⁹ Inspired by the image pyramid, an ensemble method of using different scales of features has been proposed to combine information from different scales of features to preserve different levels of image details. Initially, it was used for image classification and object detection. For example, spatial pyramid pooling²⁰ was proposed to deal with the variance in scale, size, and aspect ratio for the image classification problem. However, it was modified to detect objects with various scales, sizes, and aspect ratios. Single Shot MultiBox Detector²¹ kept six different size feature maps for object detection and achieved a better detection performance. Feature Pyramid Networks²² generated predictions at different feature levels for a single scale input image in order to take advantage of different levels of semantic information. Pyramid Scene Parsing Network²³ has been proposed for the pyramid pooling module to take advantage of prior global semantics and to capture different scales of contextual information by a parallel feature map stacking method. Deeplabv3+²⁴ used the atrous spatial pyramid pooling to replace the downsampling method to avoid the risk of potential information loss.

In this study, we present a CNN-based solution for segmenting four female pelvic organ structures from MR images both at rest and at maximal Valsalva. In the deep CNN model, a Multi-Resolution Feature Pyramid (MRFP)²⁴ module was inserted into the U-Net skip connections to capture the semantic information from different scales to improve segmentation performance in blurry regions. Depth separable convolution was used to improve the encoder network convergence. Transfer learning was applied to deal with inadequate training data. In postprocessing, a level set method was used to further improve the CNN performance. The novelty of our work could be summarized in three areas. First, it represents a novel application for pelvic organ segmentation both at rest and at maximal Valsalva in women with and without POP, based on a deep learning method with MR images. Second, it is a novel design to combine MRFP with U-Net for blurry region segmentation of medical images. We proved its effectiveness in blurry pelvic organ segmentation of high-variance MR images in POP. Third, we applied a postprocessing method to deal with the failure cases and further improve segmentation performance. As a result, compared with existing segmentation methods, our method achieves the best performance.

2. MATERIALS AND METHODS

2.A. Data population and processing

We used 24 sagittal pelvic floor MR series of 12 subjects from the Michigan Pelvic Floor Research Collection that had

been obtained with the approval of the institutional ethics review committee in case-control studies of POP. The subjects included six anterior vaginal prolapse and six posterior vaginal prolapse cases. Three women with and three women without a uterus were included per group. Supine, multiplanar MR imaging was performed in both resting and stressed states (maximal Valsalva when the patient attempts to increase the intra-abdominal pressure in order to push the pelvic organs out through the vaginal canal). All of the studies were scanned with a 3T superconducting magnet (Philips Medical Systems Inc, Bothell, WA, USA) with accompanying software (v. 2.5.1.0). In the sagittal plane, at rest, of each subject 30 slices were taken in a field of view of 200×200 mm, with a thickness of 4 mm per slice and a spacing between slices of 1 mm; at maximal Valsalva, due to the time limitation for the subjects to hold the stressed status, of each subject 14 slices were taken of scanning range 360×360 mm with a thickness of 6 mm per slice and a spacing of 1 mm.²⁵ The annotation of uterus, rectum, bladder, and levator ani muscle was accomplished based on previous anatomic work²⁶ using 3D Slicer software (v.3.4.2009-10-15). The annotation was accomplished by one expert and reviewed by another senior expert. Some preprocessing steps were applied to reduce the variance between these data. All of the slices were interpolated to the same interval in height and width dimensions. These images were then resampled into 256×256 pixel sizes for CNN model training. As there were a total of 24 sagittal pelvic floor MR series from 12 subjects and a total of 528 MR images, the different datasets were assigned as 12 3D MR series (264 images) from six subjects for training, six 3D MR series (132 images) from three subjects for validation, and six 3D MR series (132 images) from three subjects for testing. The organ occurrence rate in the training data is shown in Table I. The uterus had the lowest occurrence rate, and the bladder had the highest occurrence rate.

2.B. Convolutional neural network structure

The main conceptual framework for our CNN model is illustrated in Fig. 2. The model had an encoder-decoder network structure.^{10,27} When constructing the encoder network, we adopted the Xception²⁸ structure with residual connections. To extract different scales' context information, we used the MRFP module in the skip connections between the encoder and decoder, which will be introduced in the following subsection.

TABLE I. Organ occurrence rate in training data

Organ	Uterus	Rectum	Bladder	Levator
Number of occurrence	103	152	197	112
Number of total images	256	256	256	256
Presence rate	0.40	0.59	0.77	0.44

2.B.1. Multiresolution feature pyramid

To merge context information at multiple scales, we needed these operations to have fields of view of different sizes. Larger kernel size and dilated convolution are two options. Since the parameter quantity increases drastically as the increase of kernel sizes, we adopted dilated convolution. Each MRFP module consists of four dilated convolutional layers and one average pooling layer (Fig. 2). We used 1×1 convolution with dilation 1, 3×3 convolution with dilation 1, 3×3 convolution with dilation 2, and 3×3 convolution with dilation 3 to perceive context information at scales of 1×1 , 3×3 , 5×5 , and 7×7 . All feature maps in different branches were concatenated together for the decoder network. A convolution layer was used to mix the feature maps from different scales. Therefore, the MRFP module is capable of capturing multiscale contextual information. It was applied to all five shortcut connections in our model.

2.B.2. Encoder network structure

The encoder network (Fig. 3) is essential for feature extraction as well as for segmentation. Our encoder network adopted the Xception idea,²⁸ which takes advantage of depth separable convolution to achieve the decomposition of ordinary convolution into channelwise convolution and pointwise convolution. Customization of the model structure was proposed with modification on the downsampling operation. To

preserve more detail, we replaced the pooling layers with a convolution of stride 2. Besides, we used fewer layers in the Middle Flow to avoid overfitting.

2.C. Postprocessing method

The level set is a partial differential equation (PDE)-based method. A curve could be defined as $\phi(t,x,y)$, and after giving an initialization, the curve evolves based on image-driven forces. The PDE equation is as follows²⁹:

$$\frac{\partial \phi}{\partial t} = \nabla \phi \cdot F, \phi(0, x, y) = \phi_0 \tag{1}$$

where t is the iteration times, x and y are image coordinates, $\phi_0 = 0$ defines the initial segmentation, and F is the velocity field. To be specific, in postprocessing we used the level set method to improve the segmentation organ by organ. Using the bladder as an example, before applying the level set method, we first computed the minimum 3D boundary that includes the CNN-based bladder segmentation. This 3D boundary was then used to crop the 3D data including the bladder from the original 3D MR data. Finally, with CNN-based bladder segmentation as the initialization, we applied the level set method to the cropped MR data slice-by-slice for bladder segmentation. During model testing, compared with the ground truth, we evaluated our results using Dice Similarity Coefficient (DSC) metric and we kept the results of the level set method if they are better than the initial results. In practical applications, since ground truth values are not

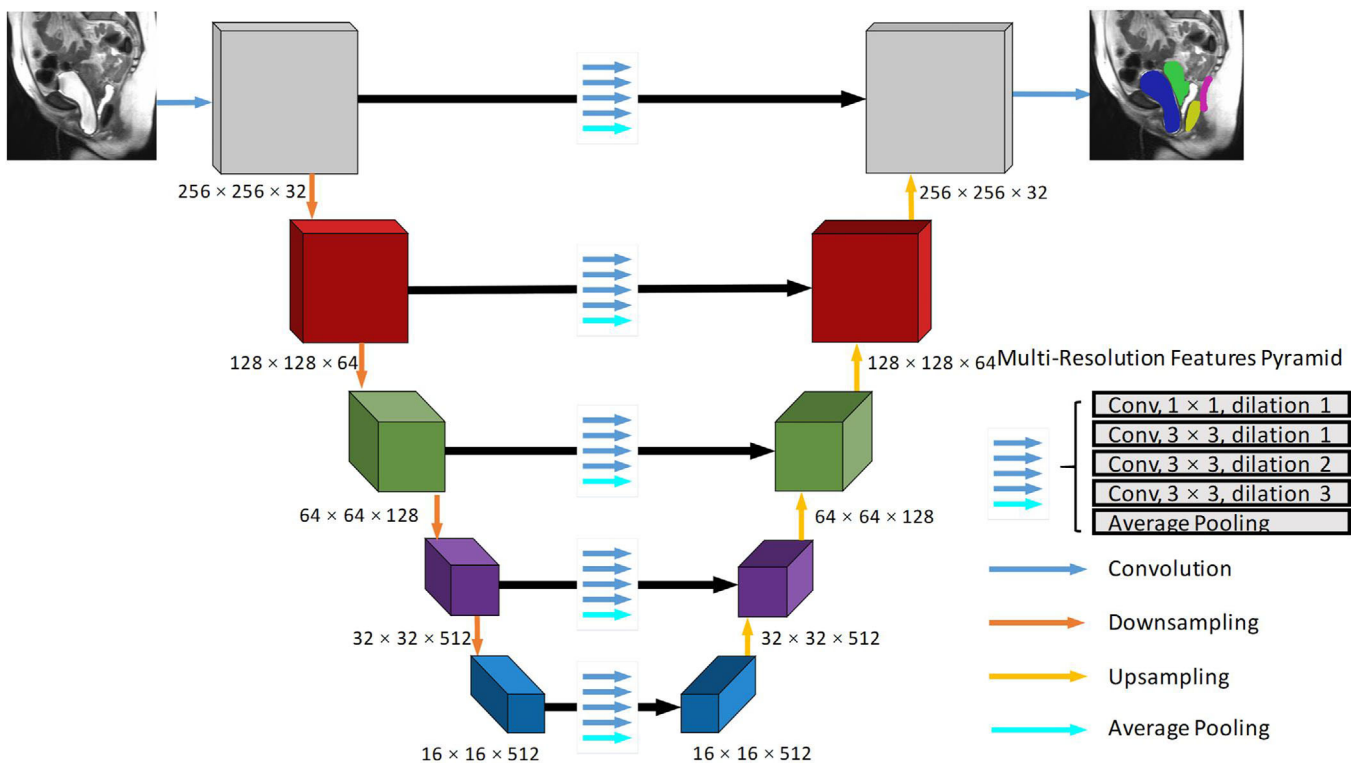


FIG. 2. Convolutional neural network model structure. Feature maps of skip connection and upsampling branches were combined using a concatenation method. [Color figure can be viewed at wileyonlinelibrary.com]

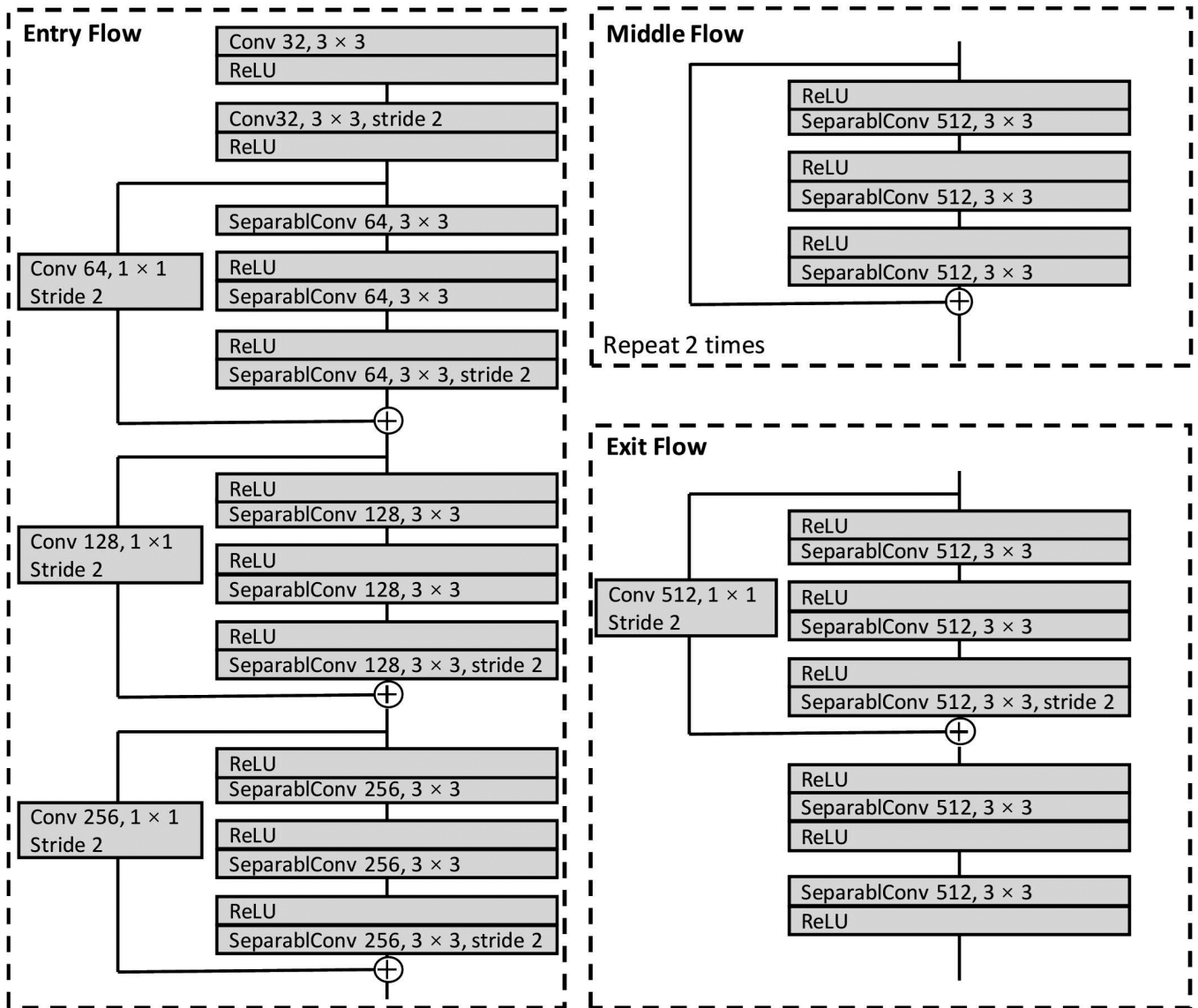


Fig. 3. Diagram illustrating the structure of the encoder network.

available, users need to determine whether the CNN model makes acceptable predictions. When users find the predictions provided by the CNN model to be unacceptable, such as the MR image segmentation is far beyond the normal range, the level set method will be applied for postprocessing, although we will only keep the better final result. For convenience, we used the *morphological_chan_vese*³⁰ function in the scikit-image library.³¹

2.D. Loss function and metrics

We investigated two different loss functions for model training, that is, pixelwise multiclass cross entropy loss (CE) and multiclass Dice loss (DL):

$$DL = 1 - 2 \frac{\sum_{l=1}^N \sum_n t_{ln} p_{ln}}{\sum_{l=1}^N \sum_n (t_{ln} + p_{ln})} \quad (2)$$

$$CE = \sum_{l=1}^N \sum_n (-t_{ln} \log(p_{ln})) \quad (3)$$

where $N = 5$ in our case, representing the background, uterus, rectum, bladder, and levator ani muscle classes, t_{ln} is the ground truth labeling on the n th pixel position for class l , and p_{ln} is the prediction result on the n th pixel position for class l .

Four metrics were used for individual organ segmentation evaluation, that is, the DSC, Average Symmetric Surface Distance (ASD), Relative Absolute Volume Difference (RAVD), and Organ Detection Recall (ODR). Following the definition of DL, the DSC is defined as follows:

$$DSC = 2 \frac{\sum_{l=1}^N \sum_n t_{ln} p_{ln}}{\sum_{l=1}^N \sum_n (t_{ln} + p_{ln})} \times 100 \quad (4)$$

And the ASD is defined as follows:

$$ASD = 2 \frac{1}{|S_T| + |S_P|} \left(\sum_{s_t \in S_T} \min_{s_p \in S_P} \|s_t - s_p\|_2 + \sum_{s_p \in S_P} \min_{s_t \in S_T} \|s_p - s_t\|_2 \right) \quad (5)$$

where S_T and S_P are the surface of the ground truth and model prediction, respectively, and s_t and s_p are corresponding points in them. The RAVD is defined as follows:

$$RAVD = \frac{|V_T - V_P|}{V_T} \times 100 \quad (6)$$

where V_T and V_P are the volume of ground truth and model prediction, respectively. The ODR is defined as follows:

$$ODR = \frac{TP}{TP + FP} \times 100 \quad (7)$$

where TP is the number of images in which an organ is correctly detected and FP is the number of images in which the same organ is not correctly detected.

2.E. Experiments

The experiment setup was summarized as below. Experiments were implemented with Keras (v.2.2.0) using Python (v.3.5.0). Adam solver was used to minimize the loss function. Our choice for the learning rate was 0.0001, with a learning rate decay of 0.98 after each epoch. A total of 800 epochs were used for training. We used an NVIDIA 1080Ti graphic card to enable the parallel computing process, with a batch size of 4. To reduce overfitting because of insufficient data, we used data augmentation. The augmentation techniques included image rotation, shear and shift, sharpening, blurring, and contrast normalization. Before images were fed to the CNN model, they were set to zero mean and unit standard variance. The Xception encoder network was trained on a cardiac structure segmentation dataset³² for transfer learning.

Experiments were conducted as follows. First, we compared DL with the CE function. Second, we compared the proposed method's performance with three other advanced segmentation methods, that is, Deeplabv3+,³³ U-Net,¹⁰ and FCN-8s.³⁴ Deeplabv3+³³ is a state-of-the-art semantic segmentation method, FCN-8s³⁴ has obtained state-of-the-art results on a PASCAL VOC 2012 Segmentation dataset, and U-Net¹⁰ is a classical biomedical segmentation method which won a challenge competition in 2015. Third, we quantified the effectiveness of the Xception encoder network and the MRFP module using ablation studies. Compared with U-Net with the Xception (U-Net-X), and U-Net with MRFP (U-Net-M), our model used U-Net with both the Xception and MRFP (U-Net-XM). Fourth, we investigated the effects of the MRFP module among different skip connections between the encoder and decoder networks. In our model, as the encoder has five downsampling stages, there are five corresponding skip

connections, which are the first to fifth skip connection from top to bottom in Fig. 2. Our model used MRFP in all the five connections so we called it U-Net-XM₁₂₃₄₅. We compared our model with its three variants, that is, U-Net-XM₁₂₃, U-Net-XM₁₃₅, and U-Net-XM₃₄₅. Finally, we used the level set method to improve the results of all segmentation methods in the second experiment.

3. RESULTS

3.A. Loss function comparison

The DL function obtained a much better segmentation result (Table II), both with and without pretraining. Hence, in the following training, we compared different methods using the DL function. The model with pretraining showed better performance than without pretraining under both loss function configurations. The pretraining improved the average DSC from 64.0 to 65.6 when using DL. However, the pretraining operation exhibited the “butterfly effect”, which means the model performance improved more in the postprocessing step (Table VII), as discussed in Section 3.E.

3.B. Performance comparison with other advanced segmentation methods

The proposed method yielded better results with respect to the DSC than the other three methods (Table III). Our model without pretraining had an average DSC of 64.0, winning in three of four individual tasks (uterus, rectum, and bladder). FCN-8s showed better performance on the rectum, but its average DSC was only 58.2. Compared with Deeplabv3+ (60.2), FCN-8s (58.2), and U-Net (54.8), our model achieved an average DSC that was 3.8, 5.8, and 9.2 points higher than them, respectively. However, our model with pretraining did not exhibit better bladder segmentation performance than the model without pretraining because the bladder of one subject was outside the normal range [Fig. 5(e)]. Segmentation of this subject was improved in the postprocessing step (see Section 3.E).

We also compared the model performances using the ODR and the RAVD (Table IV). Our model obtained the best RAVD, but did not show a distinct advantage with respect to the ODR. The ODR is the proportion of images with this organ that were correctly detected of the total number of images with this organ. The results indicate our model does not have a better organ detection ability. However, our model showed a markedly better segmentation performance (Table III), which means that for the images that were correctly detected, our model had results closer to the ground truth. A comparison of the models' predictions is shown in Fig. 4.

3.C. Ablation study

Ablation experiments were performed to quantify the effectiveness of the MRFP and the encoder network. The

TABLE II. Model performance comparison using different loss functions

Methods	Uterus		Rectum		Bladder		Levator		Average	
	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD
DL(+)	55.0 (9.3)	5.2 (1.7)	64.1 (17.6)	2.5 (1.3)	82.7 (16.5)	1.6 (0.5)	60.8 (7.4)	2.3 (1.4)	65.6	2.9
DL(*)	53.5 (18.3)	6.6 (4.6)	62.0 (17.9)	2.7 (1.1)	84.8 (10.0)	1.6 (0.5)	55.6 (9.4)	3.6 (2.5)	64.0	3.6
CE(+)	37.3 (9.8)	7.8 (2.6)	57.7 (21.6)	3.3 (1.7)	84.5 (11.7)	1.6 (0.5)	50.6 (8.3)	10.1 (11.1)	57.6	5.7
CE(*)	40.4 (19.4)	10.8 (5.6)	56.4 (16.7)	3.5 (1.4)	84.4 (10.9)	1.6 (0.3)	45.3 (13.7)	10.1 (13.6)	56.6	6.5

Best performance (mean value) will be highlighted in the use of bold. For DSC, a higher number means the better performance, while for ASD a lower number means the better performance. Units: Dice similarity coefficient (DSC) in %, and average surface distance (ASD) in mm. (+) means with pretraining, and (*) means without pretraining. Number in the () is the standard deviation.

TABLE III. Models' performance comparison with other advanced segmentation methods

Methods	Uterus		Rectum		Bladder		Levator		Average	
	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD
Proposed (+)	55.0 (9.3)	5.2 (1.7)	64.1 (17.6)	2.5 (1.3)	82.7 (16.5)	1.6 (0.5)	60.8 (7.4)	2.3 (1.4)	65.6	2.9
Proposed (*)	53.5 (18.3)	6.6 (4.6)	62.0 (17.9)	2.7 (1.1)	84.8 (10.0)	1.6 (0.5)	55.6 (9.4)	3.6 (2.5)	64.0	3.6
DeepLabv3+	45.0 (11.2)	7.3 (3.3)	58.9 (16.8)	3.0 (1.2)	83.3 (10.7)	1.9 (0.5)	53.4 (13.1)	3.8 (2.0)	60.2	4.0
FCN-8s	39.8 (14.9)	6.9 (4.8)	65.6 (11.7)	2.5 (1.0)	80.0 (13.7)	1.9 (0.7)	47.4 (16.9)	9.5 (11.6)	58.2	5.2
U-Net	45.0 (16.0)	14.3 (7.6)	42.0 (27.2)	4.7 (2.6)	77.2 (23.2)	2.9 (2.2)	54.6 (11.0)	5.2 (5.6)	54.8	6.8

Best performance (mean value) will be highlighted in the use of bold. For DSC, a higher number means the better performance, while for ASD a lower number means the better performance. Units: Dice similarity coefficient (DSC) in %, and average surface distance (ASD) in mm. (+) means with pretraining, and (*) means without pretraining. Number in the () is the standard deviation.

difference between U-Net-M and U-Net is the use of MRFP. The difference between U-Net-X and U-Net is the use of Xception encoder network. Therefore, the difference between our model (U-Net-XM₁₂₃₄₅) with U-Net-X or U-Net-M is the use of MRFP or Xception, respectively. The result is summarized in Table V.

The DSC of U-Net-M, compared with U-Net, increased from 54.8 to 61.6, an increase of 6.8 points; the DSC of our model, compared with U-Net-X, increased from 58.7 to 64.0, an increase of 5.3 points; the DSC of U-Net-X, compared with U-Net, increased from 54.8 to 58.7, an increase of 3.9 points; the DSC of our model, compared with U-Net-M, increased from 61.6 to 64.0, an increase of 2.4 points. This proved the effectiveness of MRFP when used with U-Net or

U-Net-X. Besides, MRFP made a larger contribution to the final segmentation performance. For each organ, with respect to the DSC, MRFP made a larger contribution to the uterus and the bladder than for the rectum and the levator.

3.D. Different multiresolution features pyramid module combinations comparison

The detailed segmentation results are summarized in Table VI. For the average DSC, our model (U-Net-XM₁₂₃₄₅) obtained almost the same results with U-Net-XM₃₄₅ and U-Net-XM₁₃₅, while it was 2.4 points higher than U-Net-XM₁₂₃. For individual organ segmentation, our model achieved almost the same results with U-Net-XM₃₄₅ and U-Net-XM₁₃₅ for the

TABLE IV. Models' performance comparison using other metrics

Methods	Uterus		Rectum		Bladder		Levator		Average	
	ODR	RAVD	ODR	RAVD	ODR	RAVD	ODR	RAVD	ODR	RAVD
Proposed (+)	84.8 (16.2)	34.5 (14.9)	100 (0.0)	41.0 (31.6)	91.6 (4.6)	10.8 (8.8)	95.1 (5.8)	19.9 (16.0)	92.9	26.6
Proposed (*)	84.5 (14.1)	43.3 (7.0)	94.6 (8.0)	37.6 (21.0)	98.0 (4.4)	8.6 (5.0)	91.6 (8.6)	22.2 (19.4)	92.1	28.0
DeepLabv3+	87.2 (15.7)	52.7 (35.6)	94.5 (8.0)	27.8 (16.4)	94.7 (5.4)	6.2 (5.4)	80.1 (16.2)	30.4 (16.3)	89.1	29.3
FCN-8s	84.1 (14.9)	61.9 (34.9)	96.7 (7.5)	21.0 (8.2)	98.0 (2.8)	11.9 (17.9)	91.7 (8.6)	52.8 (29.3)	92.6	36.9
U-Net	94.2 (9.9)	47.5 (33.5)	82.5 (21.6)	55.2 (26.7)	90.6 (5.7)	23.6 (24.2)	91.4 (8.9)	20.5 (16.3)	90.0	36.6

Best performance (mean value) will be highlighted in the use of bold. For ODR a higher value means the better performance and for RAVD a lower value means the better performance. Units: Organ detection recall (ODR) in %, and relative absolute volume difference (RAVD) in %. (+) means with pretraining, and (*) means without pretraining. Number in the () is the standard deviation.

TABLE V. Ablation study results

Methods	Uterus		Rectum		Bladder		Levator		Average	
	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD
Proposed (*)	53.5 (18.3)	6.6 (4.6)	62.0 (17.9)	2.7 (1.1)	84.8 (10.0)	1.6 (0.5)	55.6 (9.4)	3.6 (2.5)	64.0	3.6
U-Net-M	49.5 (11.3)	9.8 (6.8)	63.6 (16.6)	3.3 (2.2)	79.4 (18.4)	2.0 (0.8)	52.8 (10.0)	6.8 (10.0)	61.6	4.9
U-Net-X	41.2 (13.4)	11.0 (6.5)	63.4 (14.8)	2.9 (1.5)	76.1 (28.3)	2.9 (2.8)	54.2 (8.0)	3.3 (1.8)	58.7	5.0
U-Net	45.0 (16.0)	14.3 (7.6)	42.0 (27.2)	4.7 (2.6)	77.2 (23.2)	2.9 (2.2)	54.6 (11.0)	5.2 (5.6)	54.8	6.8

Best performance (mean value) will be highlighted in the use of bold. For DSC, a higher number means the better performance, while for ASD a lower number means the better performance. Units: Dice Similarity Coefficient (DSC) in %, and average surface distance (ASD) in mm. Proposed model is the U-Net-XM₁₂₃₄₅. (*) means without pretraining. Number in () is the standard deviation.

uterus and bladder, and slightly worse results for the rectum, and slightly better results for the levator. The rectum results improved using the postprocessing technique in Section 3.E (Table VII). With respect to the ASD, our model obtained the best results. Besides, U-Net-XM₁₂₃ obtained better results than U-Net-XM₃₄₅ and U-Net-XM₁₃₅.

3.E. Postprocessing improvement

We improved all CNN methods' results with the level set method. A comparison of the models' predictions is shown in Fig. 5. We demonstrated the resegmentation results by organs. Since the levator and rectum were usually connected and showed no visible edges, it was difficult to segment them using the level set method. Therefore, the uterus [Figs. 5(a) and 5(b)], rectum [Figs. 5(c) and 5(d)], and bladder [Figs. 5(e) and 5(f)] were used for comparison. With the deep learning model's prediction as prior knowledge, the level set method remedied the failure cases to a certain extent [Figs. 5(a), 5(c), and 5(e)]. However, compared with the deep learning method, the level set method did not provide better segmentation results in some general cases [Figs. 5(b), 5(d), and 5(f)] even with the deep learning model's prediction as initialization.

Final segmentation results of CNN methods after postprocessing are summarized in Table VII. Our model obtained the best DSC and ASD results for both individual organs and the overall average. The model without pretraining achieved an average DSC of 66.1 points, outperforming other methods with 4.0

to 9.7 points. Our model with pretraining obtained the highest average DSC (69.4 points) and best average ASD (2.9 mm).

4. DISCUSSION

4.A. Convolutional neural network application to pelvic organ prolapse analysis

Our work represents a novel application for female pelvic organ segmentation both at rest and at maximal Valsalva in women with and without POP, using a CNN method with MR images. In the end, we presented a hybrid solution for simultaneous uterus, rectum, bladder, and levator ani muscle segmentation and showed good results qualitatively and quantitatively. There are some differences with previous investigations.³⁵⁻⁴⁰ Different modalities of medical imaging techniques have their own advantages. Two groups used ultrasound images to accomplish levator hiatus segmentation using the fully CNN (FCN) and U-Net.^{37,41} Wang et al.³⁸ and He et al.³⁹ investigated prostate, rectum and bladder segmentation using axial view computed tomography based on a multistage FCN. Techniques including dilated convolution⁴² and full-resolution residual network⁴³ were also investigated to deal with the blurry edges of objects by capturing a larger field of view information. The level set technique as a shape prior has been considered previously for natural image segmentation.⁴⁴

Although MR imaging is the golden standard for analyzing POP, it is quite challenging, even for clinical experts, to segment pelvic organs in MR images at rest and at maximal

TABLE VI. Models' performance comparison for different MRFP configurations

Methods	Uterus		Rectum		Bladder		Levator		Average	
	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD
Proposed (*)	53.5 (18.3)	6.6 (4.6)	62.0 (17.9)	2.7 (1.1)	84.8 (10.0)	1.6 (0.5)	55.6 (9.4)	3.6 (2.5)	64.0	3.6
U-Net-XM ₁₂₃	50.9 (14.0)	8.3 (6.5)	66.1 (12.4)	2.6 (1.0)	83.7 (12.1)	2.0 (0.7)	45.8 (11.0)	4.7 (3.3)	61.6	4.4
U-Net-XM ₁₃₅	54.1 (10.6)	7.5 (5.6)	65.6 (14.5)	2.8 (1.4)	84.6 (12.5)	1.6 (0.7)	52.3 (9.6)	6.5 (4.3)	64.2	4.6
U-Net-XM ₃₄₅	53.6 (16.2)	10.8 (6.3)	65.5 (12.8)	3.5 (2.2)	84.8 (11.8)	1.6 (0.7)	52.6 (12.1)	3.3 (1.8)	64.1	4.7
U-Net-X	41.2 (13.4)	11.0 (6.5)	63.4 (14.8)	2.9 (1.5)	76.1 (28.3)	2.9 (2.8)	54.2 (8.0)	3.3 (1.8)	58.7	5.0

Best performance (mean value) will be highlighted in the use of bold. For DSC, a higher number means the better performance, while for ASD a lower number means the better performance. Units: Dice similarity coefficient (DSC) in %, average surface distance (ASD) in mm. Proposed model is the U-Net-XM₁₂₃₄₅. (*) means without pretraining. Number in () is the standard deviation.

TABLE VII. Model performance comparison after using the level set method

Methods	Uterus		Rectum		Bladder		Levator		Average	
	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD	DSC	ASD
Proposed (+)	65.3 (3.8)	5.4 (1.9)	66.3 (15.0)	2.1 (0.9)	85.6 (10.0)	1.6 (0.4)	60.8 (7.4)	2.3 (1.4)	69.4	2.9
Proposed (*)	58.3 (18.6)	6.6 (5.1)	65.8 (14.4)	2.4 (1.1)	84.8 (10.1)	1.7 (0.5)	55.6 (9.4)	3.6 (2.5)	66.1	3.6
Deeplabv3+	52.0 (14.8)	9.2 (5.5)	59.8 (17.0)	3.3 (2.0)	83.3 (10.7)	1.9 (0.5)	53.4 (13.1)	3.8 (2.0)	62.1	4.5
FCN-8s	46.0 (18.3)	8.3 (6.2)	66.0 (11.8)	2.3 (0.8)	80.7 (12.2)	1.9 (0.5)	47.4 (16.9)	9.5 (11.6)	60.0	5.6
U-Net	47.6 (15.0)	11.8 (16.3)	47.6 (22.8)	7.3 (4.4)	80.8 (15.5)	2.6 (1.2)	54.6 (11.0)	5.2 (5.6)	56.4	5.3

Best performance (mean value) will be highlighted in the use of bold. For DSC, a higher number means the better performance, while for ASD a lower number means the better performance. Units: Dice similarity coefficient (DSC) in %, and average surface distance (ASD) in mm. (+) means with pretraining, and (*) means without pretraining. Number in () is the standard deviation.

Valsalva of women with and without POP. Our deep learning model's performance is also limited by the imaging quality, the stress state, the prolapse status, and the training set size, etc. For example, the difficulty changes with segmentation from different views.³⁵ Prolapse is a downward displacement and deformation of pelvic organs, and thus its analysis is usually done from sagittal views. However, it might be more difficult for both humans and computer models to segment the uterus, levator, and rectum in the sagittal view compared with the axial view, in which the smaller organs have a higher occurrence rate. For the MR images in the sagittal view, the rest images have a thickness of 4 and 1 mm spacing. At maximal Valsalva, the stress images have a thickness of 6 and 1 mm spacing. The difficulty increases when segmenting small or thin organs, such as the levator ani and the rectum. The organs of women with POP also showed more variance than those of healthy women at maximal Valsalva compared to resting state, that is, bladders of prolapsed women might become longer at maximal Valsalva, which is very different from the bladder segmentation of men. Besides, we only included 24 sagittal MR series of 12 subjects, and images of six subjects were used for model training, limiting the deep learning model's performance. Despite these challenges, nevertheless, our deep learning model still obtained the best performance compared with other methods (Table VII).

4.B. Effectiveness analysis of different components

The effectiveness of the MRFP module is illustrated by the ablation experiments. As shown in Table V, the average DSC of U-Net-M improved by 6.8 points compared with U-Net. The average DSC of our model improved by 5.3 points compared with U-Net-X. These results suggest that MRFP is efficacious for different encoder networks. Comparing the DSC for individual organs (Table V), MRFP made larger improvements for the uterus and bladder than for the rectum and levator, because no information is obtained on the edge between the levator and rectum, as shown in Figs. 4 and 5. It is even tricky for humans to segment the rectum and levator. Models with different MRFP combinations (Table VI) revealed that our model U-Net-XM₁₂₃₄₅ had almost the same average DSC as U-Net-XM₃₄₅ and U-Net-XM₁₃₅, but a better

result on average ASD. U-Net-XM₁₂₃ achieved a lower average DSC than U-Net-XM₃₄₅ and U-Net-XM₁₃₅, but a better average ASD. A possible explanation for these observations is that MRFP on higher-order (fourth and fifth) skip connections could improve model convergence, while MRFP on lower-order (first and second) skip connections could smooth the segmentation results. In the end, our model U-Net-XM₁₂₃₄₅ achieved the best results for both average DSC and ASD, and it is therefore the recommended design.

The effectiveness of the Xception encoder network is shown in Table V. The average DSC of U-Net-X was 3.9 points higher than that of U-Net. The average DSC of our model was 2.4 points higher than that of U-Net-M on average DSC. This proved the importance of an encoder network, and a better encoder network is useful to improve segmentation.

The effectiveness of pretraining was proved in Tables III and VII. We can conclude the pretraining made a larger contribution to the uterus and levator segmentation than to the rectum and bladder segmentation. We used a cardiac MR dataset for pretraining, but a larger pelvic MR dataset might give better results. It also means more training data could be helpful to improve segmentation.

The effectiveness of the postprocessing method is shown in Tables III and VII. It also proved useful for all the CNN methods in our experiments. However, these improvements were based on using the CNN model prediction as prior knowledge. The level set method made improvements for some failure cases, such as for the examples in Figs. 5(a), 5(c), and 5(e). However, for general cases, the level set method did not provide better segmentation than the CNN method even with the CNN prediction as initialization, such as for the examples in Figs. 5(b), 5(d), and 5(f). This suggests that the CNN method has an advantage in blurry region segmentation due to training with "big data." On the contrary, since it is often challenging to collect medical imaging data and to label them, the non-deep learning method could be useful to improve the model performance to some extent. So far, whether postprocessing has improved the results needs to be compared with the ground truth. This means that it is up to the user to determine whether or not postprocessing is needed. Fortunately, comparison is a much easier task than manual segmentation. But it points to the fact that we can

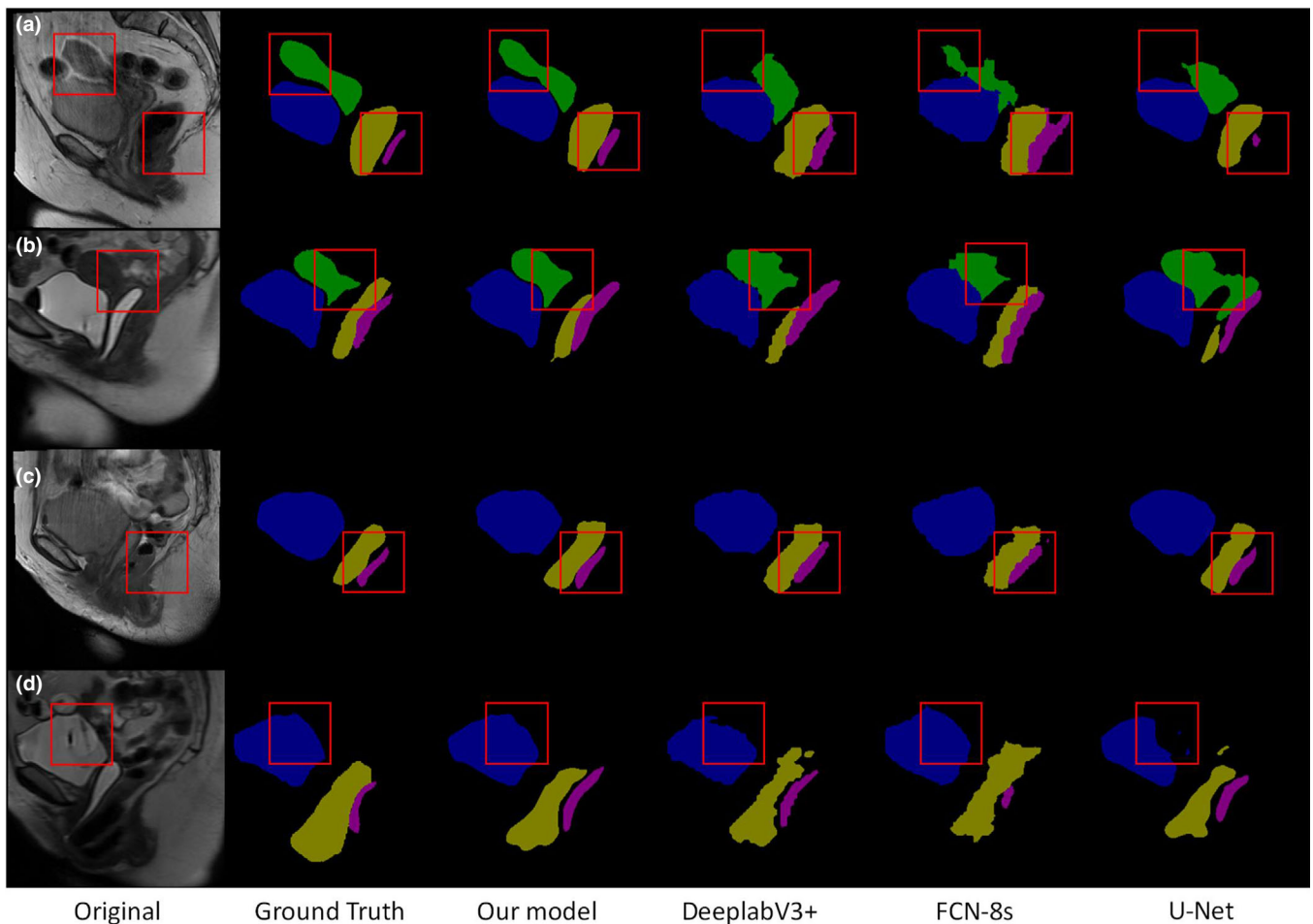


FIG. 4. A comparison of segmentation results among our model, Deeplabv3+, FCN-8s, and U-Net. (a) Resting example with uterus. (b) Stressed example with uterus. (c) Resting example without uterus. (d) Stressed example without uterus. Results of different methods were compared with the ground truth labeling. [Color figure can be viewed at wileyonlinelibrary.com]

integrate the level set method into the CNN workflow to achieve better and faster segmentation.

4.C. Segmentation performance analysis

We improved the segmentation performance from three aspects. First, we used the MRFP module to improve the blurry region segmentation on pelvic MR images. The average DSC when using MRFP increased from 54.8 to 61.6 points (Table V). Second, we built the encoder network based on the Xception idea and transfer learning technique. With the Xception, our model's performance increased from 61.6 to 64.0 points (Table V). Pretraining process improved the average DSC from 64.0 to 65.6 points (Table III). However, the pretraining operation contributed to more improvements (3.8 points) in the postprocessing step (Table VII). Third, we introduced the level set method as a postprocessing technique to deal with the limited training data and high-variance problems. Using postprocessing, our model with pretraining improved from 65.6 to 69.4 points on average DSC (Table VII). With respect to the DSC, our model outperformed other methods with 7.3 to 13.0 points. Additionally,

we compared the models' performances using the ODR and the RAVD (Table IV). Our model did not show a distinct advantage with respect to the ODR, which means our model does not detect more organs than other methods. Nevertheless, our model showed better segmentation performance (Table III), suggesting that with respect to the organs that were correctly detected, our model's results are closer to the ground truth.

The segmentation performance was ordered as follows: bladder > rectum > uterus > levator. The results of the bladder were markedly better, because the bladder has larger size, and clearer boundary than that of other organs. The rectum is easy to detect since its ODR results were higher compared to the levator and uterus (Table IV). Half of the subjects did not have a uterus, which further exacerbated the shortage of training data and the imbalance of the data, resulting in a low ODR. However, our model could predict whether there is a uterus from the subject level evaluation. After postprocessing, the highest DSC for the uterus was 65.3, which exhibited the largest improvement, as shown in Tables III and VII. The levator ani had the worst segmentation results, since it has the smallest size and does not have a clear boundary; identifying

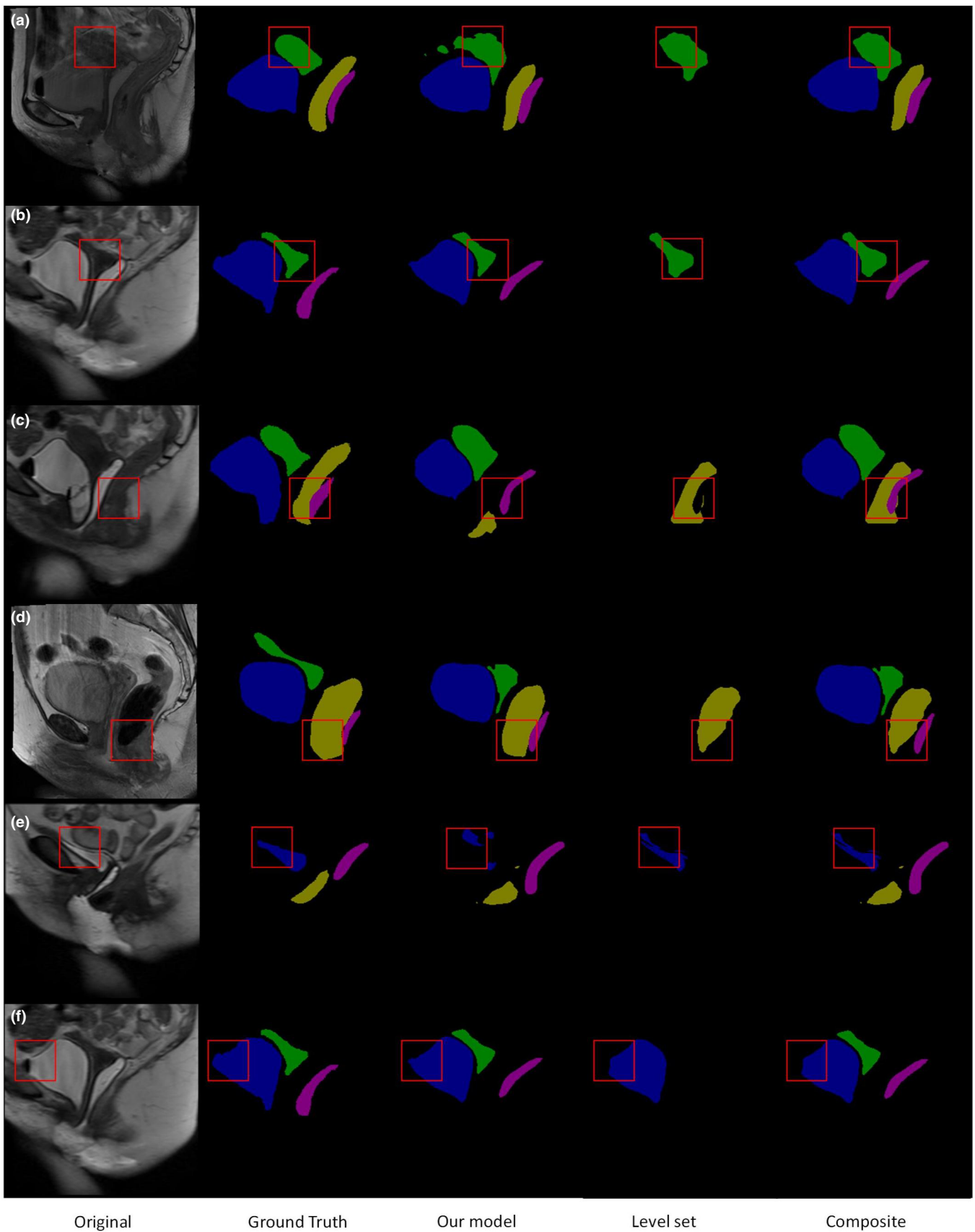


FIG. 5. Examples of re-segmentation results using the level set method. (a) and (b) Uterus re-segmentation. (c) and (d) Rectum re-segmentation, (e) and (f) Bladder re-segmentation. The composite results were obtained by replacing the models' predictions with the level set results on the corresponding organ. Results were compared with the ground truth labeling. [Color figure can be viewed at wileyonlinelibrary.com]

the levator ani is always a challenge, even for experienced clinicians.

5. CONCLUSIONS

To segment pelvic organs at rest and at maximum Valsalva (stress), we proposed a novel CNN design by integrating the MRFP module into an encoder-decoder model. This proved useful to address the blurry soft tissue boundary problem on MR images in POP. Together with the Xception encoder network and model pretraining, our model obtained better segmentation results than Deeplabv3+, FCN-8s, and U-Net. Moreover, due to the limited training data problem, a level set method was used to improve the segmentation of failure cases. Future directions include feature fusion between 2D and 3D CNNs to exploit spatial context information as discussed by Isensee et al.⁹ Model pretraining with unlabeled data using unsupervised or self-supervised methods, which could take advantage of more data, can also potentially improve the segmentation quality.

ACKNOWLEDGMENTS

This work was supported by the National Science Foundation of China General Program grant 31870942, US Public Health Service grants R01 HD038665 and P50 HD044406.

CONFLICT OF INTEREST

The authors have no relevant conflict of interest to disclose.

^{a)} Author to whom correspondence should be addressed. Electronic mail: jiajia.luo@pku.edu.cn.

REFERENCES

- Boyles SH, Weber AM, Meyn L. Procedures for pelvic organ prolapse in the United States, 1979-1997. *Am J Obstet Gynecol.* 2003;188:108-115.
- Subak LL, Waetjen LE, van den Eeden S, Thom DH, Vittinghoff E, Brown JS. Cost of pelvic organ prolapse surgery in the United States. *Obstet Gynecol.* 2001;98:646-651.
- Chen L, Ashton-Miller JA, DeLancey JOL. A 3D finite element model of anterior vaginal wall support to evaluate mechanisms underlying cystocele formation. *J Biomech.* 2009;42:1371-1377.
- Luo J, Chen L, Fenner DE, Ashton-Miller JA, DeLancey JOL. A multi compartment 3-D finite element model of rectocele and its interaction with cystocele. *J Biomech.* 2015;48:1580-1586.
- Hoyte L, Ye W, Brubaker L, et al. Segmentations of MRI images of the female pelvic floor: a study of inter- and intra-reader reliability. *J Magn Reson Imaging.* 2011;33:684-691.
- Ma Z, Jorge RNM, Mascarenhas T, Tavares JMRS. Segmentation of female pelvic cavity in axial T2-weighted MR images towards the 3D reconstruction. *Int J Num Meth Biomed Eng.* 2012;28:714-726.
- Ma Z, Tavares JMRS, Jorge RN, Mascarenhas T. A review of algorithms for medical image segmentation and their applications to the female pelvic cavity. *Comput Meth Biomech Biomed Eng.* 2010;13:235-246.
- Malone HR, Syed ON, Downes MS, D'Ambrosio AL, Quest DO, Kaiser MG. Simulation in neurosurgery: a review of computer-based simulation environments and their surgical applications. *Neurosurgery.* 2010;67:1105-1116.
- Isensee F, Jaeger PF, Full PM, Wolf I, Engelhardt S, Maier-Hein KH. Automatic cardiac disease assessment on cine-MRI via time-series segmentation and domain specific features. In: Pop M, Sermesant M, Jodoin PM, et al., eds. *Statistical Atlases and Computational Models of the Heart. ACDC and MMWHS Challenge.* Berlin: Springer International Publishing; 2016:120-129.
- Ronneberger O, Fischer P, Brox T. U-Net: convolutional networks for biomedical image segmentation. In: Navab N, Hornegger J, Wells WM, Frangi AF, eds. *Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015.* Cham: Springer International Publishing; 2015:234-241.
- Kamnitsas K, Ledig C, Newcombe VFJ, et al. Efficient multi-scale 3D CNN with fully connected CRF for accurate brain lesion segmentation. *Med Image Anal.* 2017;36:61-78.
- Havaei M, Davy A, Warde-Farley D, et al. Brain tumor segmentation with deep neural networks. *Med Image Anal.* 2017;35:18-31.
- Liskowski P, Krawiec K. Segmenting retinal blood vessels with deep neural networks. *IEEE Trans Med Imaging.* 2016;35:2369-2380.
- Lecun Y, Bengio Y. *Convolutional Networks for Images, Speech, and Time-Series.* Cambridge, MA: MIT Press; 1995.
- LeCun Y, Bengio Y, Hinton G. Deep learning. *Nature.* 2015;521:436-444.
- He K, Zhang X, Ren S, Sun J. Deep residual learning for image recognition. In: 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR); 2016:770-778.
- Milletari F, Navab N, Ahmadi S. V-Net: fully convolutional neural networks for volumetric medical image segmentation. In: 2016 fourth international conference on 3D Vision (3DV); 2016:565-571.
- Zhou Z, Rahman Siddiquee MM, Tajbakhsh N, Liang J. UNet++: a nested U-Net architecture for medical image segmentation. In: Stoyanov D, Taylor Z, Carneiro G, et al., eds. *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support.* Cham: Springer International Publishing; 2018:3-11.
- Zeiler MD, Fergus R. Visualizing and understanding convolutional networks. In: Fleet D, Pajdla T, Schiele B, Tuytelaars T, eds. *European Conference on Computer Vision - ECCV 2014.* Cham: Springer International Publishing; 2014:818-833.
- He K, Zhang X, Ren S, Sun J. Spatial pyramid pooling in deep convolutional networks for visual recognition. *IEEE Trans Pattern Anal Mach Intell.* 2015;37:1904-1916.
- Liu W, Anguelov D, Erhan D, et al. SSD: single shot multibox detector. In: Leibe B, Matas J, Sebe N, Welling M, eds. *European Conference on Computer Vision 2016.* Berlin: Springer International Publishing; 2016:21-37.
- Lin TY, Dollár P, Girshick R, He K, Hariharan B, Belongie S. Feature Pyramid Networks for Object Detection. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR); 2017:936-944.
- Zhao H, Shi J, Qi X, Wang X, Jia J. Pyramid scene parsing network. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR); 2017:6230-6239.
- Chen L, Papandreou G, Kokkinos I, Murphy K, Yuille AL. DeepLab: semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected CRFs. *IEEE Trans Pattern Anal Mach Intell.* 2018;40:834-848.
- Luo J, Larson KA, Fenner DE, Ashton-Miller JA, DeLancey JOL. Posterior vaginal prolapse shape and position changes at maximal Valsalva seen in 3-D MRI-based models. *Int Urogynecol J.* 2012;23:1301-1306.
- Luo J, Ashton-Miller JA, DeLancey JOL. A model patient: female pelvic anatomy can be viewed in diverse 3-dimensional images with a new interactive tool. *Am J Obstet Gynecol.* 2011;205:391.e1-2.
- Badrinarayanan V, Kendall A, Cipolla R. SegNet: a deep convolutional encoder-decoder architecture for image segmentation. *IEEE Trans Pattern Anal Mach Intell.* 2017;39:2481-2495.
- Chollet F. Xception: Deep Learning with Depthwise Separable Convolutions. In: 2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR); 2017:1800-1807.
- Osher S, Sethian JA. Fronts propagating with curvature-dependent speed: algorithms based on Hamilton-Jacobi formulations. *J Comput Phys.* 1988;79:12-49.
- Márquez-Neila P, Baumela L, Alvarez L. A morphological approach to curvature-based evolution of curves and surfaces. *IEEE Trans Pattern Anal Mach Intell.* 2014;36:2-17.

31. van der Walt S, Schönberger JL, Nunez-Iglesias J, et al. scikit-image: image processing in Python. *PeerJ*. 2014;2:e453.
32. Bernard O, Lalonde A, Zotti C, et al. Deep learning techniques for automatic mri cardiac multi-structures segmentation and diagnosis: is the problem solved? *IEEE Trans Med Imaging*. 2018;37:2514–2525.
33. Chen LC, Zhu Y, Papandreou G, Schroff F, Adam H. Encoder-decoder with atrous separable convolution for semantic image segmentation. In: Ferrari V, Hebert M, Sminchisescu C, Weiss Y, eds. *European Conference on Computer Vision - ECCV 2018*. Berlin: Springer International Publishing; 2018:833–851.
34. Shelhamer E, Long J, Darrell T. Fully convolutional networks for semantic segmentation. *IEEE Trans Pattern Anal Mach Intell*. 2017;39:640–651.
35. Zhou S, Nie D, Adeli E, et al. Fine-grained segmentation using hierarchical dilated neural networks. In: Frangi AF, Schnabel JA, Davatzikos C, Alberola-López C, Fichtinger G, eds. *Medical Image Computing and Computer Assisted Intervention - MICCAI 2018*. Cham: Springer International Publishing; 2018:488–496.
36. Nie D, Gao Y, Wang L, Shen D. ASDNet: attention based semi-supervised deep networks for medical image segmentation. In: Frangi AF, Schnabel JA, Davatzikos C, Alberola-López C, Fichtinger G, eds. *Medical Image Computing and Computer Assisted Intervention - MICCAI 2018*. Cham: Springer International Publishing; 2018:370–378.
37. Bonmati E, Hu Y, Sindhwani N, et al. Automatic segmentation method of pelvic floor levator hiatus in ultrasound using a self-normalizing neural network. *J Med Imaging*. 2018;5:1–8.
38. Wang S, He K, Nie D, Zhou S, Gao Y, Shen D. CT male pelvic organ segmentation using fully convolutional networks with boundary sensitive representation. *Med Image Anal*. 2019;54:168–178.
39. He K, Cao X, Shi Y, Nie D, Gao Y, Shen D. Pelvic organ segmentation using distinctive curve guided fully convolutional networks. *IEEE Trans Med Imaging*. 2019;38:585–595.
40. Nie D, Wang L, Gao Y, Lian J, Shen D. STRAINet: spatially varying stochastic residual adversarial networks for MRI pelvic organ segmentation. *IEEE Trans Neural Netw Learn Syst*. 2019;30:1552–1564.
41. Wang N, Wang Y, Wang H, Lei B, Wang T, Ni D. Auto-context fully convolutional network for levator hiatus segmentation in ultrasound images. *ISBI*; 2018:1479–1482.
42. Xia HY, Sun WF, Song SX, Mou XW. Md-Net: multi-scale dilated convolution network for CT images segmentation. *Neural Process Lett*. 2020;51:2915–2927.
43. Shah MP, Merchant SN, Awate SP. MS-Net: mixed-supervision fully-convolutional networks for full-resolution segmentation. In: Frangi AF, Schnabel JA, Davatzikos C, AlberolaLopez C, Fichtinger G, eds. *Medical Image Computing and Computer Assisted Intervention - MICCAI 2018, volume 11073 of Lecture Notes in Computer Science*. Cham: Springer; 2018:379–387.
44. Han YM, Zhang SH, Geng ZQ, Qin W, Zhi OY. Level set based shape prior and deep learning for image segmentation. *IET Image Proc*. 2020;14:183–191.